

Structural distance and evolutionary relationship of networks

Anirban Banerjee

Max Planck Institute for Molecular Genetics,
Innestrasse 63-73, 14195 Berlin, Germany

banerjee@molgen.mpg.de

December 2, 2009

Abstract

Evolutionary mechanism in a self-organized system cause some functional changes that force to adapt new conformation of the interaction pattern between the components of that system. Measuring the structural differences one can retrace the evolutionary relation between two systems. We present a method to quantify the topological distance between two networks of different sizes, finding that the architectures of the networks are more similar within the same class than the outside of their class. With 43 metabolic networks of different species, we show that the evolutionary relationship can be elucidated from the structural distances.

Author's Summary

Studying the common features and universal qualities shared by a particular class of networks in biological and other domain is one of the important aspects for evolutionary study. To measure the topological commonality, we propose a method that quantify the difference between two network structures of different sizes. Applying this measurement procedure we show that the networks from the same domain have more similarities than others. Due to the interplay between the network architecture and dynamics, biological and other networks from different areas followed by different dynamics have different structures, where networks constructed from same evolutionary process have structural similarities. We analyze 43 metabolic networks from different species and mark the prominent separation of three groups, Bacteria, Archaea and Eukarya. That is well captured in our findings that support the other cladistic results based on gene content and ribosomal RNA sequences. Thus we show that how evolutionary relationship can be elucidated from the structural distances measured by our method.

Introduction

In self-organized systems, some hidden dynamics play a role to organize the connections between the components of that system. Due to the interplay between the structure and dynamics, biological and other networks from different areas followed by different dynamics are expected to have different structures while networks constructed from the same evolutionary process have structural similarities. From structural aspects, it is important to find the answer to the question of regarding the existence of a prominent difference between different types of networks, e.g., metabolic, protein-protein interaction, power grid, co-authorship or neural networks. Studying the common features and universal qualities shared by a particular class of biological networks is one of the important aspects for evolutionary studies. In that regard, one can think about the differences between the networks within a same class, for instance among all metabolic networks, and also pose a question: are two metabolic networks from two different species, being evolutionary close more similar than others?

In the last few years different notions of the graph theory have been applied and new heuristic parameters have been introduced to analyze the network topology, for instance degree distribution, average path length, diameter, betweenness centrality, transitivity or clustering coefficient etc. (see [23] for details). Those quantities, which manage to capture particular and specific properties of the graph but not all the qualitative aspects, are not good representers of the structure and hence, with those parameters it is not possible to distinguish or compare different real networks from the point of view of topology and source of formation. Nowadays it is a fashion to categorize networks according to their degree distribution which is the distribution of k_n , the number of vertices that have degree n . It has been observed that most of the real networks have power-law degree distribution [1, 8, 13, 15, 16, 25], thus this notion also fails to distinguish networks from different systems. Hence focusing on particular and specific features is not enough to reveal the structural complexity in biological and other networks.

In this article, we propose a method to quantify the structural differences between two networks. We also show that the evolutionary relationships between the networks can be derived from their topological similarities captured by this quantification. We apply this method to the metabolic networks of 43 species and show that the phylogenetic evidences can be traced from the measurement of their structural distances.

The basic tool we use to characterize the qualitative topological properties of a network is the normalized graph Laplacian (in short Laplacian) spectra. Not only the global properties of the graph structure are reflected from the Laplacian spectrum, local structures produced by certain evolutionary processes, like motif joining or duplication are also well captured by the eigenvalues of this operator [2, 3, 4]. Distribution of the spectrum has been considered as a qualitative representation of the structure of a graph [5]. Comparative studies on real networks are difficult because of their complicated, irregular structure and

different sizes. For any graph, all eigenvalues of the graph Laplacian operator are bounded within a specific range (0 to 2). This creates the advantage to compare the spectral plots of the graphs with different sizes. Spectral plots that can distinguish the networks from different origins have been used to classify the real networks from different sources[6]. Since networks constructed from the same evolutionary process produce very similar spectral plots, the distance between spectral distributions can be considered as a measurement of the structural differences. So it can be used to study the evolutionary relation between the networks. Here, we quantify this distance with the help of an existing divergence measure (Jensen-Shannon divergence) between two distributions, what we consider as the quantitative distance measure of those two structures.

Spectrum of graph Laplacian

The normalized graph Laplacian (henceforth simply called the Laplacian) operator (Δ) has been introduced on an undirected and unweighted graph Γ , representing a network with a vertex set $V = \{i : i = 1, \dots, N\}$. For functions $v : V \rightarrow \mathbb{R}$, graph Laplacian¹ [3, 17, 18] has been defined as

$$\Delta v(i) := v(i) - \frac{1}{n_i} \sum_{j, j \sim i} v(j). \quad (1)$$

A nonzero solution u of the equation $\Delta u - \lambda u = 0$ is called an eigenfunction for the eigenvalue λ . Δ has N eigenvalues, some of them may occur with higher multiplicity. The eigenvalues of this operator are real and non-negative (because Δ is selfadjoint with respect to the product $(u, v) := \sum_i n_i u(i)v(i)$ and $(\Delta u, u) \geq 0$). The smallest eigenvalue $\lambda_0 = 0$ always, since $\Delta u = 0$, for any constant function u and the multiplicity of this eigenvalue is equal to the number of components with the graph. The highest eigenvalue λ_{N-1} is bounded above i. e. $\lambda_{N-1} \leq 2$, the equality holds *iff* the graph is bipartite². Another property of the spectra of a bipartite graph is if λ is an eigenvalue, $2 - \lambda$ is also an eigenvalue of that graph, hence the spectral plot will be symmetric about 1. The first nontrivial eigenvalue (λ_1 for connected graph) tells us how easily one graph can be cut into two different components. For the complete connected graph all nontrivial eigenvalues will be equal to $\frac{N}{N-1}$.

Along with capturing the global topological characteristics of a network, Laplacian spectrum can reveal the local structural properties. It also has the potential to describe different evolutionary mechanisms of graph formation. For instance, a single vertex $i_0 \in \Gamma$ (the simplest motif) duplication produces eigenvalue 1, which can be found with a very high multiplicity in many biological networks, with an eigenfunction u_1 that takes nonzero values at i_0 and its duplicate j_0 with $u(i_0) = 1$, $u(j_0) = -1$, and vanishes at other vertices. Duplication

¹This operator has the spectrum like the operator investigated in [10] but it has a different spectrum than the operator $Lv(i) := n_i v(i) - \sum_{j, j \sim i} v(j)$ usually studied in the graph theoretical literature as the (algebraic) graph Laplacian (see [22] for this operator).

²The distance of λ_{N-1} from 2 reflects how the graph is far from the bipartiteness.

of an edge (the motif of size two) connecting the vertices i_1 and i_2 generates the eigenvalues $\lambda_{\pm} = 1 \pm \frac{1}{\sqrt{n_{i_1}n_{i_2}}}$, and the duplication of a chain ($i_1 - i_2 - i_3$) of length 3 produces the eigenvalues $\lambda = 1, 1 \pm \sqrt{\frac{1}{n_{i_2}}(\frac{1}{n_{i_1}} + \frac{1}{n_{i_3}})}$. The duplication of these two motifs creates the eigenvalues which are close to 1 and symmetric about 1. For certain degrees of the vertices the duplication of these motifs can generate the specific eigenvalues 1 ± 0.5 and $1 \pm \sqrt{0.5}$ which are also mostly observed in the spectrum of real networks. If we join a motif Σ , which has an eigenvalue λ with an eigenfunction u_{λ} that vanishes at a vertex $i \in \Sigma$, via identifying the vertex i with any vertex of a graph Γ , the new graph will also have the same eigenvalue λ with an eigenfunction that takes the same values as u_{λ} on Σ and vanishes on other vertices. As an example, if we join a triangle that itself has an eigenvalue 1.5 to any graph, it contributes the same eigenvalue to the new graph produced by the joining process (for more details see [2, 3, 4, 5, 7]).

Jensen-Shannon divergence as a measure for the structural distance

In discrete system, Kullback-Leibler divergence measure (KL) is defined on two probability distributions p_1 and p_2 of a discrete random variable X as

$$KL(p_1, p_2) = \sum_{x \in X} p_1(x) \log \frac{p_1(x)}{p_2(x)} \quad (2)$$

Note that Kullback-Leibler (in short K-L) divergence measure is not defined when $p_2 = 0$ and $p_1 \neq 0$ for any $x \in X$. K-L divergence is not symmetric i.e. $KL(p_1, p_2) \neq KL(p_2, p_1)$ and does not satisfy the triangle inequality, hence can not be considered as a *metric*.

Jensen-Shannon divergence measure (JS) is defined on two probability distributions p_1 and p_2 as

$$JS(p_1, p_2) = \frac{1}{2}KL(p_1, p) + \frac{1}{2}KL(p_2, p); \text{ where } p = \frac{1}{2}(p_1 + p_2) \quad (3)$$

Whereas Jensen-Shannon (in short J-S) divergence is symmetric and unlike the K-L divergence measure, it does not have any problem to be defined when one of the probability measure is zero for some value of x where the other is not (for more details see [19]). Square root of J-S divergence is a metric (for details [24]).

Here we have defined the structural distance $D(\Gamma_1, \Gamma_2)$ between two different graphs Γ_1 and Γ_2 , with the spectral distribution (of graph Laplacian) f_1 and f_2 respectively, in terms of the J-S divergence measure between f_1 and f_2 :

$$D(\Gamma_1, \Gamma_2) = \sqrt{JS(f_1, f_2)} \quad (4)$$

Theoretically there exist isospectral graphs but they are relatively rare in real networks and qualitatively quite similar in most respects. For example, all complete bipartite graphs, $K_{m,n}$ (with $m + n = \text{constant}$), have the same spectrum.

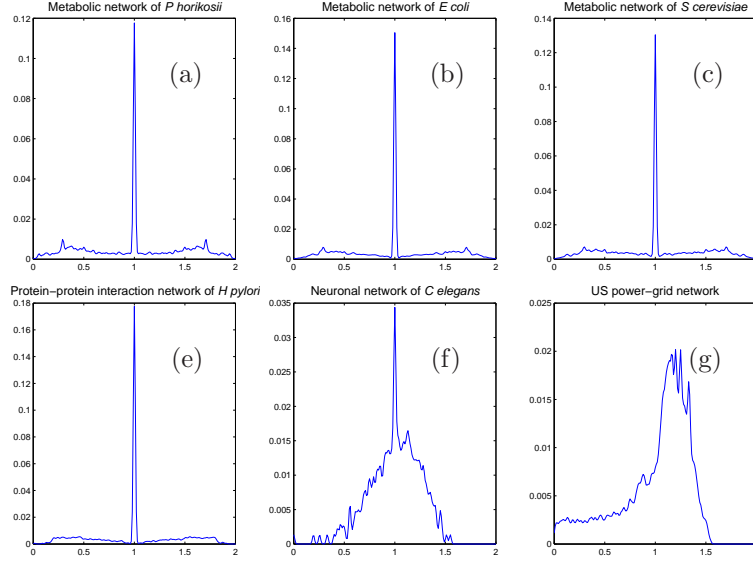


Figure 1: Spectral plots of the metabolic networks of (a) *P horikoshii*, (b) *E coli*, (c) *S cerevisiae*. The sizes of the networks are 945, 2859 and 1812 respectively. Here the nodes represent substrates, enzymes and intermediate complexes. (d) Protein-protein interaction network of *H pylori*. Network size = 710. (e) Neuronal connectivity of *C elegans*. Size of the network = 297. (f) Topology of the Western States power-grid of the United States. Network size = 4941. Here we plot the spectrum as the collection of the eigenvalues λ_i by convolving with a Gaussian kernel (with $\sigma = 0.01$). i.e. we plot $f(x) = \sum_{\lambda_i} \frac{1}{0.01\sqrt{2\pi}} \exp(-\frac{|x-\lambda_i|^2}{0.0002})$ along the vertical axis.

In this case distance between those two structure will be the same. This is one drawback of this measurement.

Results

Recalling the spectral similarities between different networks, metabolic networks are very similar to each other, and in comparison with the other networks, they are closer with the protein-protein interaction networks than the neuronal or US power-grid networks in the spectral terms [6]. Due to similar mechanisms (many metabolites or proteins have the same neighbors) of the network formation it is expected that the metabolic networks will have similar architecture with the protein-protein interaction networks rather than neuronal or power-grid networks. This phenomenon is particularly reflecting in the spectral plots (Fig.1) of the metabolic networks of *P horikoshii*, *E coli*, *S cerevisiae* with network sizes 945, 2859 and 1812 respectively, protein-protein interaction

network of *H pylori* with size 710, neuronal connectivity of *C elegans* with network size 297 and US power-grid network of size 4941 (for further reference we denote these networks by Γ_{Ph} , Γ_{Ec} , Γ_{Sc} , Γ_{Hp} , Γ_{Ce} and Γ_{PG} respectively). Now we measure the structural distances between those networks with our metric D . The differences and similarities between those networks are clearly captured by this measurement (see the Table 1). Note that each network has a different size, but nevertheless we can measure the structural distance by comparing their spectral distributions.

All the distances between these three metabolic networks are closer to each other than the protein-protein interaction network, but far from the neuronal and power-grid network. It is the same for the protein-protein interaction network. The relative distance between neuronal and power-grid networks, comparative to the other networks, is less but not as close as the one between the protein-protein interaction network and metabolic networks. These results show that we can consider our suggested metric as a suitable measure for structural differences.

Network	Γ_{Ph}	Γ_{Ec}	Γ_{Sc}	Γ_{Hp}	Γ_{Ce}	Γ_{PG}
Γ_{Ph}	0.0000	0.0904	0.0661	0.1694	0.4704	0.4704
Γ_{Ec}	0.0904	0.0000	0.0641	0.1036	0.4902	0.5074
Γ_{Sc}	0.0661	0.0641	0.0000	0.1340	0.4574	0.4738
Γ_{Hp}	0.1694	0.1036	0.1340	0.0000	0.5086	0.5380
Γ_{Ce}	0.4704	0.4902	0.4574	0.5086	0.0000	0.2429
Γ_{PG}	0.4780	0.5074	0.4738	0.5380	0.2429	0.0000

Table 1: Distance table between metabolic networks of *P horikoshii* (Γ_{Ph}), *E coli* (Γ_{Ec}), *S cerevisiae* (Γ_{Sc}); protein-protein interaction network of *H pylori* (Γ_{Hp}); neuronal connectivity network of *C elegans* (Γ_{Ce}) and US power-grid network (Γ_{PG}). All the distances are computed using the metric $D(\Gamma_1, \Gamma_2)$.

Evolutionary relationship from the distance measure

Networks constructed from the same evolutionary process are structurally close to each other. Thus, the architectures of the networks that share the same evolutionary path are expected to be more similar than others. So to a large extent, one can elucidate the evolutionary relationships between the networks within the same system from their structural distances. To verify this conviction we evolve a graph along a tree (see Fig. 2(a)) and predict the evolutionary relations among the graphs of a generation. Here we choose the initial graph A0, a scale-free network constructed by the Barabási–Albert’s model [8] ($m_0 = 5$ and $m = 3$). After a certain number of edge-rewiring, while keeping the degree of each node the same, we produce a graph of the next generation. Note that here all the graphs have not only the same degree distribution but also the same degree

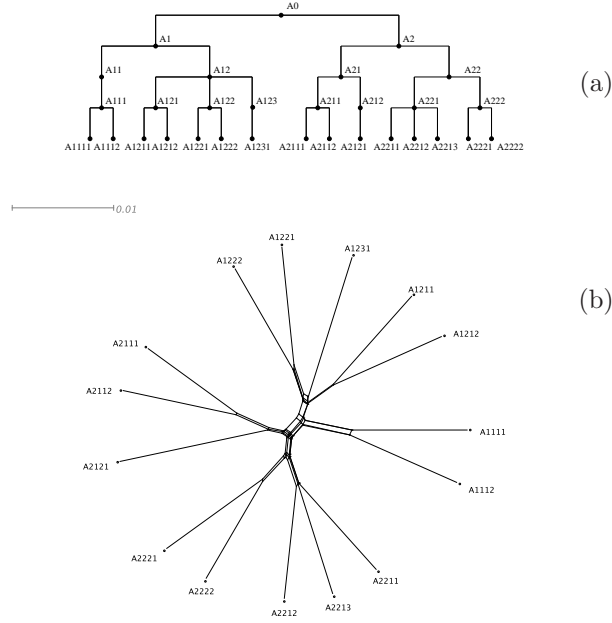


Figure 2: (a) Evolution of a graph A_0 along a definite tree: A_1 and A_2 have been produced independently in the 2nd generation with a certain evolutionary process from A_0 . In the same way, A_{11} and A_{12} have been produced from A_1 and A_{21} , A_{22} from A_2 and so on. Continuing in the same fashion, we end up with the graphs $A_{1111}, \dots, A_{2222}$ in the 5th generation. (b) The splits network for the structural distances (calculating by our proposed metric) of the graphs from the 5th generation. Each band of parallel edges indicate a split. For example, two lines represent the split $\{A_{1111}, A_{1112}\}$ versus the other graphs. This tree-like splits network shows that the evolutionary relationships among those graphs is clearly captured by our distance measure. The figure has been produced by using Neighbor-Net [9].

sequence. One can also choose any other evolutionary mechanism. But that would not make any significant difference in the result. We take all the graphs having been produced in the same generation (here we choose generation 5) and estimate the structural distances between them using our measure D (in 4). Now for these distances we produce a splits network [14], which can extract phylogenetic signals that are missed in other tree-representation. This tree-like network (see Fig. 2(b)) shows that the distances contain a prominent phylogenetic signal and clearly demonstrates the evolutionary relationships between those graphs.

Comparison with the other structural difference measures

Other methods can also be used to quantify the structural similarities of the networks. A common way to compare two graph structures is to collate the independent heuristic parameters defined on them. For this purpose, we choose the following parameters: transitivity, diameter, radius, average path length, average edge-betweenness centrality, and average node-betweenness centrality for this purpose. Now we construct a vector V_{Γ}^{para} , using the values of the parameters mentioned above from a graph Γ as the components and compute the structural difference D^{para} between two graphs Γ_1 and Γ_2 as

$$D^{para}(\Gamma_1, \Gamma_2) = \| V_{\Gamma_1}^{para} - V_{\Gamma_2}^{para} \| \quad (5)$$

The other measure D^{motif} , we consider, is based on the normalized Z score [21] of the motif of size 3 and 4. It has been shown that the networks can be categorized in different superfamily [20] based on the characteristic distribution of the relative frequency of their motifs. In the similar way, we construct a vector V_{Γ}^{motif} from a graph Γ with the values of the normalized Z score of the motif of size 3 and 4 as the components and compute the structural difference between two graphs Γ_1 and Γ_2 as

$$D^{motif}(\Gamma_1, \Gamma_2) = \| V_{\Gamma_1}^{motif} - V_{\Gamma_2}^{motif} \| \quad (6)$$

Now we compare the efficiency of the measure D with D^{motif} and D^{para} to predict the evolutionary relationships among the graphs. Like previous way we compute the matrix with the distances estimated by a particular measure mentioned above between the graphs that are produced in the 5th generation of the graph evolution along the tree (Fig. 2(a)). We use *symmetric difference*, defined by Robinson-Foulds [26], (in short R-F distance) between the tree constructed from a distance matrix using neighbor-joining method and the true tree shown in Fig. 2(a). The R-F distance between two trees is the number of bipartitions that can be found in one tree but not in other one. Since our true tree contains two internal nodes (A12 and A221) of degree 4, the neighbor joining (in short N-J) tree with all the internal nodes have degree 3 always has two bipartitions which are never present in the true tree. A N-J tree that resembles the true tree most will have a R-F distance of 2 to the true tree. Fig. 3(a), which shows three frequency distributions of such R-F distances for every measures, clearly demonstrate that the measure D is more accurate than the other two. The limited accuracy can be explained by the stochasticity in the process of graph evolution. In order to address whether the accuracy is also influenced by systematic effects, we investigate the trend in the R-F distances of the trees that are constructed using the sum of k distance matrices produced by using a particular measure over k realizations of graph evolution from the true tree. The R-F distance decreases and assumes its minimum value 2 with increasing k (Fig. 3(b)). For this particular graph evolution, the evolutionary relationships can be perfectly recovered from the information of the D -measure, if the input size become large enough. However evidently, the spectral distribution captures

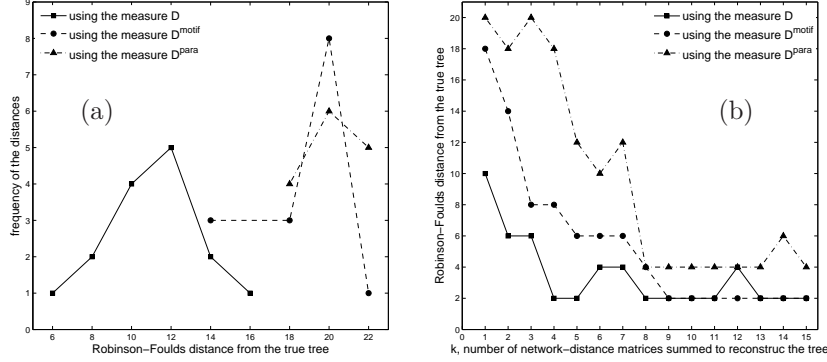


Figure 3: The measure D is more accurate than D^{motif} and D^{para} . (a) Frequency distributions of the Robinson-Foulds distances of the trees that are constructed from graph structural-distances using D , D^{motif} , D^{para} from the true tree (in Fig. 2(a)). (b) Here we plot the Robinson-Foulds (R-F) distances along the vertical axis. We produce the graph distance matrices using D , D^{motif} , D^{para} for every k realization of graph evolution. Then we sum all the k distance matrices for each measures and compute the R-F distances of the trees reconstructed from these summed matrices from the true tree.

more qualitative properties of a network than the heuristic parametric values and the expression of the small motifs do.

Evolutionary relationships between metabolic networks of 43 species

Now using our structural difference measure D we estimate the distances between the metabolic networks of 43 species and construct a distance matrix between them. Fig. 4, which is a splits network for these distances, supports that the data contained in that matrix has a substantial amount of phylogenetic signal and some parts of the data are tree-like. Due to the non-uniform evolutionary rate of topological change, to analyze the structural similarities among the networks of all those species we construct an unrooted tree from the mentioned distance matrix by using the neighbor-joining method. This tree, which resembles highly the phylogenetic tree of those 43 species, shows different clusters according to the structural similarities of the metabolic networks (see Fig. 5). The prominent separation of three groups, Bacteria, Archaea and Eukarya³. That is well captured in our findings that support the other cladistic results based on gene content [27] and ribosomal RNA sequences [28]. This is a strong evidence how evolutionary relationship is reflected from the structural similarities which are clearly captured by the measure of the spectral distances by our metric D .

³Only Yeast belongs to the group of Bacteria.

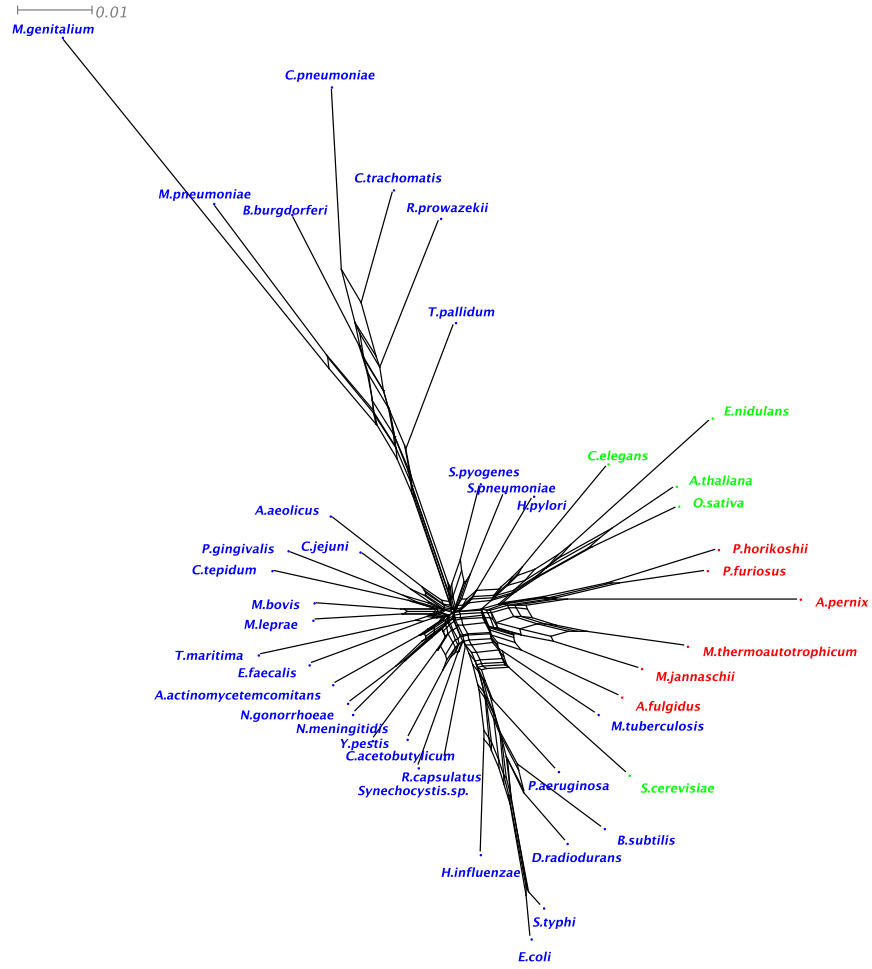


Figure 4: The splits network for the structural distances (calculating by the metric D) between the metabolic networks (of 43 species). This network shows that the distance-data is tree-like and has some phylogenetic signal. The colors, blue, green and red indicate Bacterium, Eukaryote and Archae respectively. We use Neighbor-Net [9] to produce this figure.

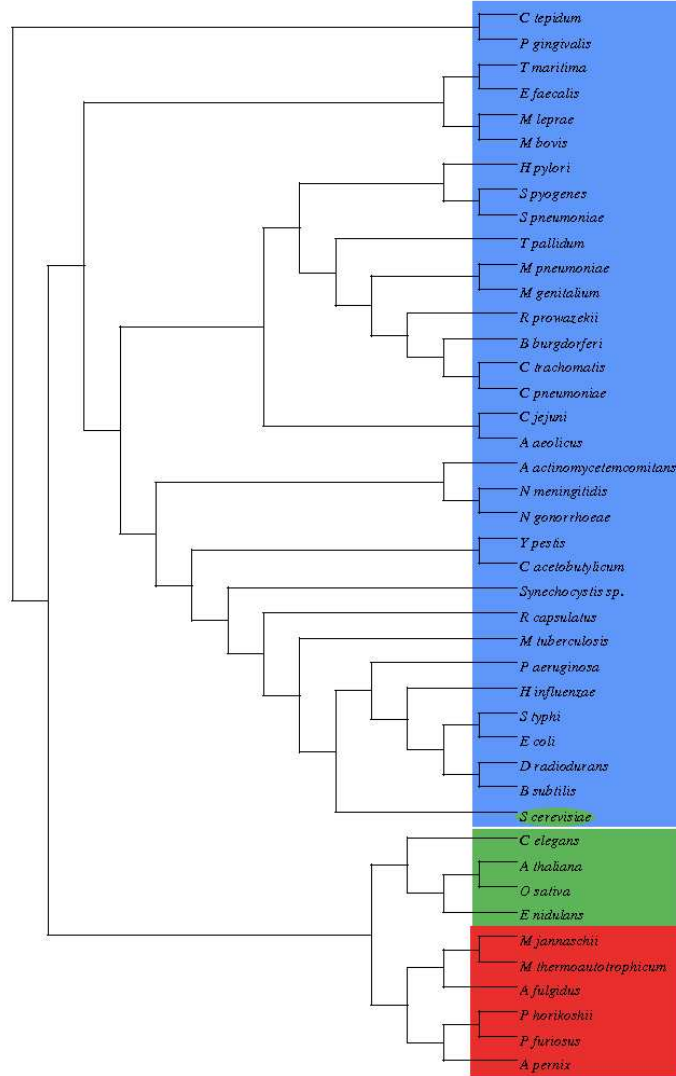


Figure 5: The un-rooted tree of metabolic networks (of 43 species) constructed with their structural distances (calculating by our proposed metric) using the neighbor-joining method. Bacterium, Eukaryote and Archae are showed by the color, blue, green and red respectively and all of them form separate cluster within the tree. Only *S. cerevisiae* belongs to a different group, Bacterium.

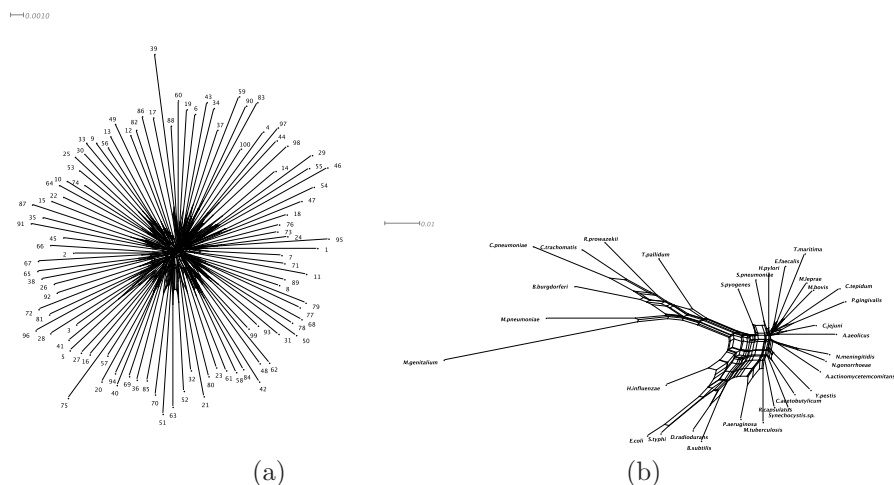


Figure 6: The splits network of the structural distances between (a) 100 networks constructed by randomly deleting 5 percent of the reactions from the metabolic network of *E. coli* and (b) metabolic networks of 32 bacteria. The star-like structure of the splits network in (a), which is very different from the splits network of bacteria in (b), shows that the data of the distance matrix merely has a phylogenetic signal and the metabolic networks of bacteria are not constructed only by mapping from the *E.coli*. We have used Neighbor-Net [9] to construct both the splits networks.

Cross validation of the tree construction against the effect of the enzyme mapping from *E. coli*

All the metabolic pathways in *E. coli* have been constructed independently in wetlab. But it is not always the case for the other bacteria. If an enzyme-specific gene that also exists in *E. coli* has been detected, the same metabolic reactions catalyzed by that enzyme are incorporated into the database. If there are no different genes which have been reported from every other bacteria and that can make significant change in the network structure, all other metabolic networks will be very similar and the detection of the phylogenetic relationship can be an artifact. In order to verify this fact, we reconstruct 100 networks by randomly deleting 5 percent of the reactions from the metabolic network of *E. coli* and produce a splits network of the distances between those 100 networks. The star-like structure of this splits network, which is very different from the splits network constructed from the structural distances between the metabolic networks of 32 bacteria, shows that the distances of those 100 networks merely have a phylogenetic signal (Fig. 6). Hence the evolutionary relationships can not be detected if all other metabolic networks are only mapped from the network of *E. coli*.

Discussion

Here we suggest a method to compare the architecture of the networks with different sizes, an aspect causing the main problem for the comparison. With a defined metric, we quantify their structural similarities based on the spectral distribution which captures the qualitative properties of the underlying graph topology which can emerge from the evolutionary process like motif duplication or joining, random rewiring, random edge deletion etc. In spite of the network reconstruction error (see source of the data), this method elucidate the evolutionary relationships between the metabolic networks constructed from 43 different species. To explore the evolutionary relationships in other domains like language and society structure and in other biological areas, this approach can also be used.

Methods

Sources of the data

In this article we use the data set which are freely available. We access the metabolic data (used in [16]) of 43 species from <http://www.nd.edu/~networks/>. At the time of database construction genomes of 25 species (18 bacteria, 2 eukaryotes and 5 archaea) had been completely sequenced while the remaining 18 species underwent this process partially. But the analysis of the errors [16] suggest that there would not be a drastic change in the final result. We use the network data for the protein-protein interaction of *Helicobacter pylori* from <http://www.cosinproject.org/> and neuronal connectivity (used in [29, 30]) of *C elegans* from <http://cdg.columbia.edu/cdg/datasets>.

Network construction from the data set

Due to incomplete sequencing of the genome of different species, many biological data are incomplete and they contain statistical errors. To capture a more appropriate (i.e. with less error) network architecture we focus on the giant component. It is very probable that this part of the network is constructed from the mostly studied metabolic pathways, hence consists more complete data and capture most of the qualitative properties of the original complete network. Moreover, in our analysis we consider the underlying undirected graphs of the real networks which are directed in many cases. The reduced graph itself carries a lot of structural information that is quite informative about the network, but one can easily extend this method to directed networks for having more accurate results.

Compute the distribution of the spectrum

After computing the spectrum of a network we convolve with a kernel $g(x, \lambda)$ and get the distribution by normalizing the function

$$f(x) = \int g(x, \lambda) \sum_k \delta(\lambda, \lambda_k) d\lambda = \sum_k g(x, \lambda_k) \quad (7)$$

Here we use the Gaussian kernel $\frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-m_x)^2}{2\sigma^2})$ with $\sigma = .01$ for all computation. Choosing other types of kernels does not change the result significantly.

Clustering of the metabolic networks by constructing an unrooted tree

Since we are interested only to get the clusters among all those metabolic networks according to their structural distance, an unrooted tree is our interest, thus the neighbor-joining method is adequate to choose for the construction. We calculate the $D(\Gamma_i, \Gamma_j)$ for each pair of those networks (Γ_i, Γ_j) and build a distance matrix. We use the software package PHYLIP [12] and SplitsTree [14] for the tree construction. The branching distance is not important for our purpose, hence we ignore the branch length while plotting the tree.

Compute the normalized Z score of a motif

The normalized Z score of a motif of a network is the normalized relative frequency of that motif, compared to its expression in the randomized version of the same network. The statistical significance of a motif σ is presented by its Z score,

$$Z_\sigma = \frac{N_\sigma^{real} - \langle N_\sigma^{rand} \rangle}{SD(N_\sigma^{rand})}, \quad (8)$$

where N_σ^{real} is the number of times the motif σ appears in the network, and $\langle N_\sigma^{rand} \rangle$ and $SD(N_\sigma^{rand})$ are the mean and standard deviation of its appearance in the ensemble of randomized networks. Hence the normalized Z score of a motif σ is $Z_\sigma / (\sum_\sigma Z_\sigma^2)^{1/2}$. Here, with the help of the software mfinder1.2, which is freely available on <http://www.weizmann.ac.il/mcb/UriAlon/>, we calculate the Z score of each motif of size 3 and 4, and normalize them over all.

Acknowledgments

The author is thankful to Martin Vingron, Thomas Manke, Roman Brinzanik, Sitabhra Sinha, Monojit Choudhur for valuable discussions. A special thank to Hannes Luz for giving the useful suggestions regarding phylogenetic tree construction. The author is also thankful to Antje Glück for the helpful comments

on preparing the manuscript. Thanks to the VolkswagenStiftung for the funding to support this project.

References

- [1] R. Albert, H. Jeong, and A. L. Barabási. Internet - diameter of the world-wide web. *Nature*, 401(6749):130131, 1999.
- [2] A. Banerjee, J. Jost. Laplacian spectrum and protein-protein interaction networks. Preprint. E-print available: arXiv:0705.3373.
- [3] A. Banerjee, J. Jost. On the spectrum of the normalized graph Laplacian. *Linear Algebra and its Applications*, 428, 3015-3022, 2008.
- [4] A. Banerjee, J. Jost. Graph spectra as a systematic tool in computational biology. *Discrete Applied Mathematics*, 157(10), 2425-2431, 2009.
- [5] A. Banerjee, J. Jost. Spectral plots and the representation and interpretation of biological data. *Theory in Biosciences*, 126(1), 15-21, 2007.
- [6] A. Banerjee, J. Jost. Spectral plot properties: Towards a qualitative classification of networks. *NHM*, 3(2), 395-411, 2008.
- [7] A. Banerjee, J. Jost. Spectral characterization of network structures and dynamics. In: N. Ganguly et al.(eds.), *Dynamics On and Of Complex Networks; Modeling and Simulation in Science, Engineering and Technology*, 117-132, Springer Birkhäuser Boston, 2009.
- [8] A. L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509512, 1999.
- [9] D. Bryant and V. Moulton. Neighbor-net: an agglomerative method for the construction of phylogenetic networks. *Mol Biol Evol* 21:255265, 2004.
- [10] F.Chung, *Spectral graph theory*, AMS, 1997
- [11] P. Erdős, A. Rényi, On random graphs. *Publ. Math. Debrecen*, 6:290-297, 1959.
- [12] J. Felsenstein. Inferring phylogenies from protein sequences by parsimony, distance, and likelihood methods. *Methods Enzymol*, 266:418-427, 1996.
- [13] R. Guimera, S. Mossa, A. Turtshi, and L. A. N. Amaral. The worldwide air transportation network: Anomalous centrality, community structure, and cities global roles. *Proc. Natl. Acad. Sci. USA*, 102(22):77947799, 2005.
- [14] D. H. Huson, *SplitsTree: analyzing and visualizing evolutionary data*. *Bioinformatics*, 14(1):68-73, 1998.

- [15] H. Jeong, S. P. Mason, A. L. Barabási, and Z. N. Oltvai. Lethality and centrality in protein networks. *Nature*, 411(6833):4142, 2001.
- [16] H. Jeong, B. Tombor, R. Albert, Z. N. Oltvai, and A. L. Barabási. The large-scale organization of metabolic networks. *Nature*, 407(6804):651654, 2000.
- [17] J. Jost, Dynamical networks in: J.F.Feng, J.Jost, M.P.Qian (eds.), *Networks: from biology to theory*, pp.35–62, Springer, 2007
- [18] J. Jost, M. P. Joy, Spectral properties and synchronization in coupled map lattices. *Phys.Rev.E* 65, 16201-16209, 2001
- [19] J. Lin, Divergence measures based on the Shanon entropy. *IEEE Trans. on Information Theory*, 37(1):145-151, January 1991.
- [20] R. Milo et al., Superfamilies of Evolved and Designed Networks. *Science* 303: 1538-1542, 2004.
- [21] R. Milo et al., Network motifs: Simple building blocks of complex networks. *Science* 298: 824827, 2002.
- [22] B. Mohar, The Laplacian spectrum of graphs. In: Y. Alavi, G. Chartrand, O. R. Oellermann, A. J. Schwenk (eds.), *Graph Theory, Combinatorics, and Applications*, pp. 871-898, Vol. 2, Ed. Wiley, 1991.
- [23] M. E. J. Newman, The structure and function of complex networks. *SIAM Review*, 45(2):167-256, 2003.
- [24] F. Österreicher and I. Vajda, A new class of metric divergences on probability spaces and its statistical applications. *Ann. Inst. Statist. Math.*, 55:639653, 2003.
- [25] S. Redner, How popular is your paper? an empirical study of the citation distribution. *European Physical Journal B*, 4(2):131134, 1998.
- [26] D. F. Robinson, L. R. Foulds, Comparison of phylogenetic trees. *Math. Biosciences*, 53: 131147, 1981.
- [27] B. Snel, P. Bork, and M.A. Huynen, Genome phylogeny based on gene content. *Nature Genet.* 21:108110, 1999.
- [28] C.R. Woese, O. Kandler, and M.L. Wheelis, Towards a natural system of organisms: proposal for the domains Archaea, Bacteria, and Eucarya. *Proc. Natl. Acad. Sci. USA* 87:45764579, 1990.
- [29] D. J. Watts and S. H. Strogatz, Collective dynamics of small-world networks, *Nature*, 393:440442, 1998.
- [30] J. G. White et al., The structure of the nervous-system of the nematode *Caenorhabditis-Elegans*, *Phil. Trans. Royal Soc. of London Series B-Bio. Sc.*, 314:1340, 1986.